

66/12/60



ASSISTANT COMMISSIONER FOR PATENTS
Washington, D.C. 20231

Case Docket No. ST9-99-037

55.09

September 21, 1999

Express Mail Label No. EL414495639US

Dear Sir:

Transmitted herewith for filing is the patent application of

Inventor(s): Mark Anthony Cesare, Tom Robert Christopher, Julie Ann Jerves, Richard Henry Mandel III

For: **METHOD, SYSTEM, PROGRAM AND DATA STRUCTURE FOR CLEANING A DATABASE TABLE**

Enclosed are:

- ☒ 34 pages of Application for Patent Including 22 Pages Specification; 1 Page Abstract; and 46 Claims
- ☒ 10 No. of Sheets of Drawings Sheet(s) of drawings (☒ informal)
- ☒ An assignment of the invention to International Business Machines Corporation. (☐ Will follow.)
- ☐ An associate power of attorney.
- ☐ A verified statement to establish small entity status under 37 CFR 1.9 and 1.27.
- ☒ Declaration and Power of Attorney. (☐ Will follow.)
- ☐ Certified copy of Patent Application No. filed from which priority is claimed under 35 U.S.C. §119.
- ☐ IDS enclosed. ☐ with references.

10/18 U.S. PTO

09/39694

09/21/99

CALCULATION OF FEES

ITEM	NO. OF CLAIMS FILED MINUS BASE*	NO. OF CLAIMS OVER BASE	X SM/LG ENTITY FEE	\$ AMOUNT	\$ FEE
A TOTAL CLAIMS FEE	46 - 20* =	26	X \$9 or X \$18	\$468	
B INDEPENDENT CLAIMS FEE**	4 - 3* =		X \$39 or X 78	\$78	
C SUBTOTAL - ADDITIONAL CLAIMS FEE (ADD FINAL COLUMN IN LINES A + B)					\$548
D MULTIPLE-DEPENDENT CLAIMS FEE			SMALL ENTITY FEE = \$130 LARGE ENTITY FEE = \$260		\$0
E BASIC FEE*			SMALL ENTITY FEE = \$380 LARGE ENTITY FEE = \$760		\$760
F TOTAL FILING FEE (ADD TOTALS FOR LINES C, D, AND E)					\$1306
G ASSIGNMENT RECORDING FEE				\$ 40	\$40
**LIST INDEPENDENT CLAIMS 1, 14, 27, 40					

Please charge my Deposit Account No. 50-0585 in the amount of \$

A copy of this sheet is enclosed.

- ☒ A check in the amount of \$ 1306 to cover the filing fee is enclosed.
- ☒ Check for \$ 40 covering the Recordation of Assignment fee enclosed.
- ☒ The Commissioner is hereby authorized to charge payment of the following fees associated with this communication or credit any overpayment to Deposit Account No. 50-0585. **A copy of this sheet is enclosed.**
- ☒ Any additional filing fees required under 37 CFR 1.16.
- ☒ Any patent application processing fees under 37 CFR 1.17.
- ☐ The Commissioner is hereby authorized to charge payment of the following fees during the pendency of this application or credit any overpayment to Deposit Account No. 50-0585. **A copy of this sheet is enclosed.**
- ☐ Any patent application processing fees under 37 CFR 1.17.
- ☐ The issue fee set in 37 CFR 1.18 at or before mailing of the Notice of Allowance, pursuant to 37 CFR 1.311(b).
- ☐ Any filing fees under 37 CFR 1.16 for presentation of extra claims.

Respectfully submitted,

David W. Victor
David W. Victor
Registration No. 39,867

Direct All Correspondence to:

David W. Victor
KONRAD, RAYNES & VICTOR, LLP
1180 S. Beverly Drive; Suite 501
Los Angeles, CA 90035

Direct Telephone Calls to:

(310) 556-7983

METHOD, SYSTEM, PROGRAM, AND DATA STRUCTURE
FOR CLEANING A DATABASE TABLE

CROSS-REFERENCE TO RELATED APPLICATIONS

5 This application is related to the following co-pending and commonly-assigned patent applications, all of which are filed on the same date herewith, and which are incorporated herein by reference in their entirety:

"Method, System, Program, And Data Structure for Transforming Database
Tables," to Mark A. Cesare, Tom R. Christopher, Julie A. Jerves, Richard H.
10 Mandel III, and having attorney docket number ST9-99-034;
"Method, System, Program, And Data Structure for Pivoting Columns in a
Database Table," to Mark A. Cesare, Julie A. Jerves, and Richard H. Mandel III,
and having attorney docket number ST9-99-035;
"Method, System, and Program for Inverting Columns in a Database Table," to
15 Mark A. Cesare, Julie A. Jerves, and Richard H. Mandel III, and having attorney
docket no. ST9-99-038; and
"Method, System, Program, And Data Structure For Cleaning a Database Table
Using a Look-up Table," Mark A. Cesare, Julie A. Jerves, and Richard H. Mandel
III, and having attorney docket no. ST9-99-036.

20

BACKGROUND OF THE INVENTION

1. Field of the Invention

The present invention relates to a method, system, program, and data structure for
25 cleaning a database table and, in particular, for performing clean operations on columns
in the database table.

2. Description of the Related Art

Data records in a computer database are maintained in tables, which are a collection of rows all having the same columns. Each column maintains information on a particular type of data for the data records which comprise the rows. A data warehouse is
5 a large scale database including millions or billions of records defining business or other types of transactions or activities. Data warehouses contain a wide variety of data that present a coherent picture of business or organizational conditions over time. Various data analysis and mining tools are provided with the data warehouse to allow users to effectively analyze, manage and access large-scale databases to support management
10 decision making. Data mining is the process of extracting valid and previously unknown information from large databases and using it to make crucial business decisions. In many real-world domains such as marketing analysis, financial analysis, fraud detection, etc, information extraction requires the cooperative use of several data mining operations and techniques.

15 Once the desired database tables have been selected and the data to be mined has been identified, transformations on the data may be necessary. Transformations vary from conversions of one type of data to another, e.g., converting nominal values into numeric ones so that they can be processed by a neural network, to definition of new attributes, i.e., derived attributes. New attributes are defined either by applying
20 mathematical or logical operators on the values of one or more database attributes. The transformed data is stored in a target database where it may then be mined using one or more techniques to extract the desired type of information necessary to make the organizational decisions. Further details of data mining are described in the International Business Machines Corporation (IBM) publication entitled "White Paper: Data Mining
25 Solutions" (IBM Copyright, 1996)

Data transformation refers to the process of filtering, merging, decoding, and translating source data to create validated data for the data warehouse and data mining tools. For example, a numeric regional code might be replaced with the name of the

093094 093094

region. Data transformations and cleansing is used when data is inconsistent or incompatible between sources. In such case, some level of data cleansing is needed to ensure data consistency and accuracy. Some of the current techniques for transforming and cleansing data include the use of an SQL WHERE clause to limit the rows extracted
5 from the source table. Further, formulas and expressions specified in the column definition window and constants and tokens are used to eliminate and modify data.

Previous versions of IBM Visual Warehouse included programs to allow users to perform numerous functions on the source data. For instance, if one database table has revenue data in U.S. dollars and another data table stores revenue data in foreign currency
10 denominations, then the foreign revenue data must be cleansed before both sets of data can be analyzed together. Transformation operations may be performed using application programs external to the database program that process and transform tables of data records. Further details of data warehousing and data transforms, are described in the IBM publications "Managing Visual Warehouse, Version 3.1," IBM document no. GC26-
15 8822-01 (IBM Copyright, January, 1998), which is incorporated herein by reference in its entirety.

Notwithstanding current programs for cleansing data, there is a need in the art to provide users greater control over operations to clean input data.

20

SUMMARY OF THE PREFERRED EMBODIMENTS

To overcome the limitations in the prior art described above, preferred embodiments disclose a method, system, program, and data structure for performing a clean operation on an input table. The input table to clean is indicated in an input data table name. At least one rule definition is processed to clean the input table. Each rule
25 definition indicates a find criteria, a replacement value, and an input data column in the input table. For each rule definition, the input data column is searched for any fields that match the find criteria. The replacement value for the particular rule definition is inserted in the fields in the input data column that match the find criteria. Subsequent applications

of additional rule definitions applied to the same input data column operate on replacement values inserted in the input data column during previously applied rule definitions.

In further embodiments, each rule definition is associated with one rule table including the find criteria and replacement value. In such case, a rule table column parameter is provided for each rule definition indicating the columns in the rule table including the find criteria and replacement value for that rule definition. In certain embodiments, two rule definitions may have the same rule table. In such case, the rule table column parameters indicate different columns in the same rule table including the find criteria and replacement value for each rule definition. In still further embodiments, a separate rule table may include the find criteria and replacement value for different rule definitions.

Still further, a rule definition may include multiple find criteria and a corresponding replacement value for each find criteria. In such case, the step of searching the input data column comprises applying each of the multiple find criteria to one field until a match occurs or none of the multiple find criteria are found to match the field content. When a match is found, the replacement value corresponding to the find criteria is inserted in the field having the matching content.

In preferred embodiments, the rule definition may define a find and replace rule, a discretization rule or a numeric clip rule. Different rule definitions may define different rule types.

In preferred embodiments, the rule definitions may be communicated from one computer system, such as a client, to a computer system including the input data table, such as a database server. The rule definitions are then executed against the input table on the database server including the input tables.

Preferred embodiments provide a data command structure including one or more rule definitions for performing different operations on the data in an input data table. The preferred embodiments provide a command structure that accommodates multiple types

of clean operations to be performed on an input data table before the input data table is written to the output table. Further, preferred embodiments allow a client to transfer clean commands including to the database server including the database for execution on the database server. This reduces network traffic as the database tables subject to the
5 clean operation do not have to be transferred between the database server and the client constructing the clean commands. Further, in preferred embodiments, the rule definitions are maintained in rule tables in the server. This further reduces network traffic as the clean command need only specify the location of rules to apply and does not have to provide tables of rules.

10

BRIEF DESCRIPTION OF THE DRAWINGS

Referring now to the drawings in which like reference numbers represent corresponding parts throughout:

FIG. 1 illustrates a computing environment in which preferred embodiments are
15 implemented;

FIG. 2 illustrates the parameters used in a transform command to clean input tables in accordance with preferred embodiments of the present invention;

FIGs. 3a, 3b, 4, and 5 illustrate examples of a rule table to clean data in accordance with preferred embodiments of the present invention;

20 FIGs. 6a, 6b, 6c, 6d, and 6e illustrate logic to clean an input data table in accordance with preferred embodiments of the present invention;

FIG. 7 illustrates an example of an input data table; and

FIGs. 8a, 8b, 8c, and 8d illustrate examples of rule table to apply to clean columns in the input data table in FIG. 7 in accordance with preferred embodiments of the present
25 invention.

667260743660

DETAILED DESCRIPTION OF THE PREFERRED EMBODIMENTS

In the following description, reference is made to the accompanying drawings which form a part hereof and which illustrate several embodiments of the present invention. It is understood that other embodiments may be utilized and structural and
5 operational changes may be made without departing from the scope of the present invention.

Computing Environment

FIG. 1 illustrates a computing environment 2 in which preferred embodiments are
10 implemented. The environment 2 includes a server 4 and client 6. The server 4 and client 6 would include an operating system, such as MICROSOFT WINDOWS 98 and WINDOWS NT, AIX, OS/390, OS/400, OS/2, and SUN SOLARIS,** and may be comprised of any suitable server and client architecture known in the art. The server 4 and client 6 include a database program 8a and 8b, wherein 8a comprises the server 4 side
15 of the database program and 8b comprises the client 6 side. The server 4 and client 6 may communicate via any communication means known in the art, such as a telephone line, dedicated cable or network line, etc, using any protocol known in the art including TCP/IP network (e.g., an Intranet, the Internet), LAN, Ethernet, WAN, System Area Network (SAN), Token Ring, etc. Alternatively, there may be separate and different
20 networks between the servers 4 and client 6.

The client/server database programs 8a, b, may be comprised of any client/server database program known in the art, such as DB2, Oracle Corporation's ORACLE 8, Microsoft SQL Server,** etc. The database programs 8a and 8b are used to access operations and perform operations with respect to information maintained in one or more
25 databases 10. The database(s) 10 would consist of multiple tables having rows and columns of data, e.g., tables 14 and 18. Further details of the architecture and operation of a database program are described in the IBM publications "DB2 for OS/390: Administration Guide, Version 5" IBM document no. SC26-8957-01 (Copyright IBM.

Corp., June, 1997) and "A Complete Guide to DB2 Universal Database," by Don Chamberlin (1998), which publications are incorporated herein by reference in its entirety.

In preferred embodiments, the clean transform program is implemented using the IBM stored procedure database program structure. A stored procedure is a block of procedural constructs and embedded SQL statements, i.e., an application program, that is stored in a database and can be called by name. Stored procedures allow an application program to execute in two parts. One part runs on the client and the other on the server. This allows one client call to produce several accesses of the database from the application program executing on the system, i.e., server including the database. Stored procedures are particularly useful to process a large number of database records, e.g., millions to billions of records, without having to transfer data between the server and client. The client stored procedure passes input information to the server stored procedure which then, executing within the database program including the database, processes numerous database records according to such client input information. The server stored procedure program is initiated by the client, and during execution the client cannot communicate with the stored procedure executing in the server. Further details of stored procedures are described in the publication "A Complete Guide to DB2 Universal Database," "A Complete Guide to DB2 Universal Database," which was incorporated by reference above.

The clean transform of the preferred embodiments is implemented as a stored procedure application program 12 in the server 4. The clean transform stored procedure 12 receives as input a name of an input table 14 in the database 10, and transform rules 16 from the client 6 specifying the clean operations to perform on the data in the named input table 14. The results of the clean operations performed by the clean transform stored procedure 12 in response to the transform rules 16 are generated into the output table 18. Alternatively, the transformed, i.e., cleaned input table is written to the database 10 to overwrite the previous version of the input table.

3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025
1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079
1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133
1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187
1188
1189
1190
1191
1192
1193
1194
1195
1196
1197
1198
1199
1200
1201
1202
1203
1204
1205
1206
1207
1208
1209
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241
1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295
1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349
1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403
1404
1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457
1458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511
1512
1513
1514
1515
1516
1517
1518
1519
1520
1521
1522
1523
1524
1525
1526
1527
1528
1529
1530
1531
1532
1533
1534
1535
1536
1537
1538
1539
1540
1541
1542
1543
1544
1545
1546
1547
1548
1549
1550
1551
1552
1553
1554
1555
1556
1557
1558
1559
1560
1561
1562
1563
1564
1565
1566
1567
1568
1569
1570
1571
1572
1573
1574
1575
1576
1577
1578
1579
1580
1581
1582
1583
1584
1585
1586
1587
1588
1589
1590
1591
1592
1593
1594
1595
1596
1597
1598
1599
1600
1601
1602
1603
1604
1605
1606
1607
1608
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619
1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649
1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1660
1661
1662
1663
1664
1665
1666
1667
1668
1669
1670
1671
1672
1673
1674
1675
1676
1677
1678
1679
1680
1681
1682
1683
1684
1685
1686
1687
1688
1689
1690
1691
1692
1693
1694
1695
1696
1697
1698
1699
1700
1701
1702
1703
1704
1705
1706
1707
1708
1709
1710
1711
1712
1713
1714
1715
1716
1717
1718
1719
1720
1721
1722
1723
1724
1725
1726
1727
1728
1729
1730
1731
1732
1733
1734
1735
1736
1737
1738
1739
1740
1741
1742
1743
1744
1745
1746
1747
1748
1749
1750
1751
1752
1753
1754
1755
1756
1757
1758
1759
1760
1761
1762
1763
1764
1765
1766
1767
1768
1769
1770
1771
1772
1773
1774
1775
1776
1777
1778
1779
1780
1781
1782
1783
1784
1785
1786
1787
1788
1789
1790
1791
1792
1793
1794
1795
1796
1797
1798
1799
1800
1801
1802
1803
1804
1805
1806
1807
1808
1809
1810
1811
1812
1813
1814
1815
1816
1817
1818
1819
1820
1821
1822
1823
1824
1825
1826
1827
1828
1829
1830
1831
1832
1833
1834
1835
1836
1837
1838
1839
1840
1841
1842
1843
1844
1845
1846
1847
1848
1849
1850
1851
1852
1853
1854
1855
1856
1857
1858
1859
1860
1861
1862
1863
1864
1865
1866
1867
1868
1869
1870
1871
1872
1873
1874
1875
1876
1877
1878
1879
1880
1881
1882
1883
1884
1885
1886
1887
1888
1889
1890
1891
1892
1893
1894
1895
1896
1897
1898
1899
1900
1901
1902
1903
1904
1905
1906
1907
1908
1909
1910
1911
1912
1913
1914
1915
1916
1917
1918
1919
1920
1921
1922
1923
1924
1925
1926
1927
1928
1929
1930
1931
1932
1933
1934
1935
1936
1937
1938
1939
1940
1941
1942
1943
1944
1945
1946
1947
1948
1949
1950
1951
1952
1953
1954
1955
1956
1957
1958
1959
1960
1961
1962
1963
1964
1965
1966
1967
1968
1969
1970
1971
1972
1973
1974
1975
1976
1977
1978
1979
1980
1981
1982
1983
1984
1985
1986
1987
1988
1989
1990
1991
1992
1993
1994
1995
1996
1997
1998
1999
2000
2001
2002
2003
2004
2005
2006
2007
2008
2009
2010
2011
2012
2013
2014
2015
2016
2017
2018
2019
2020
2021
2022
2023
2024
2025
2026
2027
2028
2029
2030
2031
2032
2033
2034
2035
2036
2037
2038
2039
2040
2041
2042
2043
2044
2045
2046
2047
2048
2049
2050
2051
2052
2053
2054
2055
2056
2057
2058
2059
2060
2061
2062
2063
2064
2065
2066
2067
2068
2069
2070
2071
2072
2073
2074
2075
2076
2077
2078
2079
2080
2081
2082
2083
2084
2085
2086
2087
2088
2089
2090
2091
2092
2093
2094
2095
2096
2097
2098
2099
2100
2101
2102
2103
2104
2105
2106
2107
2108
2109
2110
2111
2112
2113
2114
2115
2116
2117
2118
2119
2120
2121
2122
2123
2124
2125
2126
2127
2128
2129
2130
2131
2132
2133
2134
2135
2136
2137
2138
2139
2140
2141
2142
2143
2144
2145
2146
2147
2148
2149
2150
2151
2152
2153
2154
2155
2156
2157
2158
2159
2160
2161
2162
2163
2164
2165
2166
2167
2168
2169
2170
2171
2172
2173
2174
2175
2176
2177
2178
2179
2180
2181
2182
2183
2184
2185
2186
2187
2188
2189
2190
2191
2192
2193
2194
2195
2196
2197
2198
2199
2200
2201
2202
2203
2204
2205
2206
2207
2208
2209
2210
2211
2212
2213
2214
2215
2216
2217
2218
2219
2220
2221
2222
2223
2224
2225
22

58 appears in the log file and is used to identify the transform operations for which the log entry was made. The replacement rule definition(s) 60 identifies an input column from the input table 14 and an output column in the output table 18, and one or more rules to use when processing the input columns. The specified field of parameters 50
5 defines the transform rules 16 the client side of the clean transform stored procedure 20 presents to the clean transform stored procedure 12 in the server 4.

The replacement rule definitions 60 parameter comprises one or more definitions, such that each definition 80 includes sub-parameters 82-100. Thus, the rule definitions parameter 60 may specify multiple distinct rule definitions 80, each describing particular
10 operations to perform on specified columns in the input table. Below is a description of the sub-parameters included in each definition 80 a user may specify in the rule definitions parameter 60:

Input Data Column Name 82 - defines an existing column in the table having the Input Data Table Name 52 that contains the data to be operated on in accordance
15 with the rules specified in the rule table and other parameters. This parameter is required in the definition 80. The input data column must be capable of being modified, so that the clean transform stored procedure 12 will not update a value which has a constraint, such as for a unique key column or a referential constraint.

20 Output Data Column Name 84: defines the name of a column in the table having the Output Data Table Name 54 where cleaned data is placed. This parameter is optional, and the default is the input data column name 82. Thus, if no output data column or output table is specified, clean data is placed back into the copy of the input table in memory being processed.

25 Rule Table Name 86: Provides the name of the table containing the "find and replace" rule patterns. This sub-parameter is optional; if not provided the values for the column are copied to the output table, and are only modified if the option

Case 1:20-cv-00000-00000

to compress or remove white space is specified. In preferred embodiments, the rule tables are maintained in the server 4 and are specified in the rule definitions.

Rule Table Sort-Key Column Name 88: When a rule table includes multiple rules to apply to the input data column, this parameter 88 indicates a sort key-column in the rule table that provides an ordering in which the rules are applied to rows (fields) in the input data column.

Rule 90: This sub-parameter is required if the rule table name 86 is specified. This sub-parameter defines the type of rule included in the rule table to use when processing the values in the input data column. Further details of rule 90 are described below.

Rule Table Columns 92: This sub-parameter is specified if the rule table name is specified. The number of columns depend on the rules, examples of which are discussed below. If a rule table has different columns for different rule definitions, i.e., one rule table stores multiple rule definitions, then the rule table columns 92 parameter indicates those columns in the rule table including the specific rules for one rule definition.

Row Clean Indicator 94: Boolean YES/NO value. If YES is specified and there is a match in the input data column to the "find" condition, then the input value is not copied to the output table and the located matching rows are deleted from the temporary input table. Note that the setting this parameter to YES will affect subsequent attempts to process the transitional input data columns as the matching entries in the input data columns subject to the row clean indicator 94 are empty. If the value for this sub-parameter is NO, the matching entry in the input data column is copied to the output data column or the input data column. Further, any replacement value in a rule table is ignored and need not be specified when the row clean indicator is used. This sub-parameter is optional, and is ignored unless a rule table name 86 is specified.

Rule Escape Character 96: Certain characters have special meaning, such as the percent sign (%), which represents any string of zero or more characters, or the underscore character (_), which represents a single character. The use of the rule escape character in the search string followed by the special character, means that the query looks for the actual special character following the escape character, and ignores the special meaning.

White Space Indicator 98: Optional boolean value that indicates whether white spaces are ignored when searching

Numeric Tolerance 100: Specifies a tolerance value when doing a "find" operation with respect to numeric values. Thus, numbers within the "tolerance" range of the searched upon number will produce a match.

A rule table may include columns for different rule definitions. In preferred embodiments the rule table is maintained in the database program 8a for direct access by the clean transform 12. In this way, users can utilize predefined clean rules in a rule table in the database program 8a. The rule table may include one find column used in two rule definitions and different columns for the replacement values for the two rule definitions having the same find value. Thus, both rule definitions search on the same criteria, but provide different replacement values. Alternatively, the rule table may have different columns for the find criteria for two rule definitions, but the same replacement value for the different columns including the find criteria. Still further, the rule table may include multiple search criteria for a single rule definition and a corresponding search value for each of the search criteria. In this way, multiple search criteria could be applied to the same input column as part of one rule definition. Upon finding the first find value matching the field content, the replacement value corresponding to that matching find value would be inserted into the matching field.

Moreover, multiple rules in different columns of the same rule table or in different rule tables may sequentially operate against the same input data column. If more than one set of rules is to be applied to an input data column, a definition 80 must be specified for

each rule. The order of the definitions listed in the replacement rule definitions 60
parameter indicates the order in which the rules are applied to the input data column. In
this way, values in the input data column may be modified in sequence and processed
according to clean operations specified in different rule definitions. In preferred
5 embodiments, the processing of subsequent rule definitions on the same column will
depend on any replacement values inserted in the field during the application of a
previous rule definition, not the original content in the input data column. Thus, rule
definitions are applied to the copy of the input data table in memory as it is being
processed. The find operation with respect to a field will apply to any previously inserted
10 replacement value inserted into the field. If this inserted replacement value matches the
search criteria of the subsequent rule definition, then another replacement value will be
inserted into the field.

The rule sub-parameter 90 indicates the type of rule indicated in the table
identified in the rule table name sub-parameter 86. Possible values for rule type could
15 include find and replace, discretize, numeric clip or any other find and replace type
operations known in the art. A find and replace locates a field in the input data column
matching the find value. This find value is specified in one column of the rule table. The
column including the find rule is provided in the first rule table column identified in sub-
parameter 92. This rule table column would further identify a second column in the rule
20 table including a replacement value. If the find operation produces a match on the find
value, then the replacement value is inserted in the field having the content matching the
find value. The specification of the value to find in the rule table column must match the
type of data in the input data column, e.g., only numbers are allowed in numeric "find"
columns, pattern strings are allowed for character columns. Further for either numeric or
25 character data, the database null value can be used as a find or replacement value.

FIGs. 3a and 3b illustrate a possible format for rule tables for a find and
replacement rule type. FIG. 3a illustrates a two column find and replace rule table for an
input data column including character values. The find pattern is what is searched for in

601260"4636E50

the input data column and the replacement text value is what is inserted into the field matching the find pattern. FIG. 3b illustrates a similar two column find and replace rule table for an input column having numerical values.

If the rule 90 indicates a discretization type, then there is both an upper and lower bound for the find value. The find operation can specify to search for matching values in the input data column that are between the upper and lower bounds, outside of the upper and lower bounds, less than the lower bound, or greater than the upper bounds. Another column in the rule table would specify a replacement operation to perform. If the upper and lower boundaries specified character strings, then SQL rules would be used to determine whether strings in the input data column satisfy the search criteria. FIG. 4 illustrates a table format for discretization for character and numeric data types. The replacement value is inserted into every field in the input data column having a value between the upper and lower bounds. Placement of the NULL indicator in the upper or lower bounds can affect whether the find seeks all fields in the input data column greater or less than one of the bounds. The discussion below with respect to FIG. 6b and 6c explains how a NULL value in one of the find values affects the search criteria.

If the rule 90 indicates a numeric clip type, then the find operation finds fields in the input data column that are less than, equal to or greater than specified values. FIG. 5 shows that discretization includes an upper and lower values for both the find and replace. Any values in the input data column less than the lower bound and greater than the upper bound are replaced with the lower and upper replacement value, respectively. As with discretization, the use of the NULL value in one of the find values can affect the find operation, as discussed below with respect to FIG. 6d and 6e.

For find and replacement, discretization, and numeric clip rule types, the database NULL value can be used as the find or replacement value. The NULL value indicates the absence of information. The database NULL value is only allowed for a pattern find value or replacement value if the database input column allows for NULL values, i.e., the input data column was not defined with a NOT NULL clause. The database NULL value may

also be used to determine how to consider lower and upper bound values in the discretization and numeric clip operations, as described below.

FIGs. 6a, b, c, d, and e illustrate logic implemented in the clean transform stored procedure 12 (clean transform 12) to clean input data columns according to the clean parameters 50 including definitions 80 in the replacement rule definition(s) 60 parameter. These parameters 50 and sub-parameters 80 comprise the transform rules 16. As discussed, a user at the client 6 would specify certain clean operations to perform and the client side of the clean transform 20 would generate an API function call including the clean parameters 50 and definitions 80 from information the user entered in the GUI.

Control begins at block 200 with the clean transform 12 receiving the clean parameters 50 and a set of sub-parameters 80 for each rule definition in an API function call. The clean transform 12 accesses (at block 202) the input table 14 from the database 10 indicated in the input data table name 52. For each rule definition 80 specified in the replacement rule definitions 60 parameter, the clean transform 12 executes a loop to perform operations in FIGs. 6a, b, c, beginning at block 204.

In processing each rule definition, the clean transform 12 accesses (at block 206) the input data column from the input table 14 indicated in the input data column name 82 parameter for the rule definition and accesses (at block 208) the rule table in the server 4 indicated in the rule table name parameter 86 and input columns in the rule table indicated in the rule table columns parameter 92 for the rule definition. The clean transform 12 determines (at block 210) whether the rule table sort-key column name is non-empty. If so, the clean transform 12 sorts (at block 210) the rules in the rule columns for the rule definition according to the order specified in the sort-key column in the rule table. As discussed, a sort-key column may be provided if a rule table includes multiple rules to apply for the rule definition. Otherwise, if the sort key-column parameter 88 is empty, the rule columns are sorted in ascending order. After block 212 or the yes branch of block 210, the clean transform 12 determines (at block 214) whether the rule parameter 90 indicates a find and replace clean operation.

If the clean operation for the current definition is a find and replace, the clean transform 12 determines (at block 216) one or more find values or patterns (if the data type is a character) from the first column in the rule table indicated in the rule table columns parameter 92. The clean transform 12 inserts (at block 218) the replacement value/pattern in the column in the rule table having the same column number as the second column in the rule table columns 92 sub-parameter if one of the find values or patterns matches the field content. Note that if there are multiple find values/patterns in the rule table, then there is a different corresponding replacement value for each separate find/value pattern. The replacement value corresponding to the first matching find value is inserted in the field. The clean transform then returns (at block 220) to block 204 to process the next rule definition 80 in the rule definition parameters 60.

During the find and replace searching, and the searching operations for the discretization and numeric clip rule types, other sub-parameters are considered. If a rule escape character is indicated in sub-parameter 96, the clean transform 12 searches for a character matching the character following the escape character, which may be a character that usually has specific search meaning, such as a percent (%), underscore (_), comma (,) or semicolon (;). Further, if the row clean indicator sub-parameter 94 is set to YES, then any field, i.e., row, in the input data column matching the find value/pattern is not copied to the output table in the current rule definition being applied, or any further rule definitions that may apply to the field or row. If a numeric tolerance sub-parameter 100 is provided, than number fields in the input data column will return a match if the field value is equal to the find value within a range specified by the numeric tolerance. Further, if the ignore white space indicator sub-parameter 98 is YES, then white spaces are removed from the field when determining whether the field matches the find value/parameter; otherwise, white spaces are considered characters.

If the rule indicated in the rule sub-parameter 90 is discretization (at block 240), then the clean transform 12 determines (at block 242) the upper bound, lower bound and replacement values from the accessed rule table. The upper, lower, and replacement

columns to access from the rule table maintained in the database program 8a are indicated in the first, second, and third columns in the rule table columns sub-parameter 92, respectively. The clean transform 12 determines (at block 244) whether the upper and lower bounds are both NULL. If so, the clean transform 12 does the replacement (at
5 block 246) in every field in the current input data column. Otherwise, the clean transform 12 begins a loop at block 248 to perform for every field in the input data column. In this loop, the clean transform 12 executes a nested loop between blocks 249 and 274 to consider each rule in the rule table, when there are multiple rules. For each rule, the clean transform 12 determines (at block 250) whether the accessed field in the input data
10 column is NULL. If so, the clean transform 12 determines (at block 252) whether either the upper or lower bounds are NULL. If so, the clean transform 12 inserts (at block 254) the replacement value in the rule table into the field; otherwise, no replacement value is inserted (at block 256). From blocks 254 or 256, control transfers to block 258 where the clean transform 12 performs another iteration of the loop beginning at block 248 for the
15 next field (row) in the input data column.

If the field in the input data columns is not NULL (from the NO branch at block 250), then the clean transform 12 determines (at block 260) whether both the upper and lower bounds are NULL. If so, the clean transform 12 inserts the replacement value into the field and proceeds to block 258 to consider any further rows in the input data column.
20 Otherwise, the clean transform 12 determines (at block 264) whether only the upper bound is NULL. If so, the clean transform 12 inserts (at block 266) the replacement value into the current field if the field in the input data column is greater than the lower bound minus the numeric tolerance for numeric data types where a numeric tolerance is specified. If no numeric tolerance is provided in parameter 100 or the input column is a
25 character type, then there is no consideration of numeric tolerance when performing the find operation. If only the lower bound is NULL (at block 268), then the clean transform 12 inserts (at block 270) the replacement value into the current field if the field in the input data column is less than the upper bound plus any numeric tolerance for numeric

2025 RELEASE UNDER E.O. 14176

data types where a numeric tolerance is specified. If neither the upper nor lower bound are NULL, then the clean transform 12 inserts (at block 272) the replacement value in the field if the field value is less than or equal to the upper bound plus any provided numeric tolerance AND less than or equal to the lower bound minus any provided numeric

5 tolerance. After applying one rule in the rule table at blocks 254, 256, 262, 266, 270 or 272, the clean operation 12 then considers any further rules at block 274 for the current row in the input data table. After finding a match or considering all rules in the rule table for a given field, the clean transform proceeds (at block 258) to consider the next row (field) in the input data column.

10 After all rows in the input data column are considered from blocks 246 or 258, the clean transform proceeds (at block 276) to block 204 to execute the next rule definition against a specified input data column.

If the rule indicated in the rule parameter 90 is a numeric clip (at block 300 in FIG. 6d), then the clean transform determines (at block 302) the upper bound, lower

15 bound, and replacement value in the columns of the rule table indicated in the rule table columns parameter 92. The clean transform then determines (at block 302) whether both upper and lower bounds are NULL. If so, the clean transform 12 skips (at block 304) the input data column, and inserts no replacement values into any row. Otherwise, the clean transform 12 begins a loop at block 306 to process each row in the input data column.

20 The clean transform 12 begins a nested loop at block 307 to consider each rule in the rule table, if there are multiple rules. The clean transform 12 determines (at block 308) whether the field content is NULL. If so, the clean transform 12 does not insert (at block 310) the replacement value for any into the field and proceeds (at block 312) to consider the next row in the column until all rows are processed. If only the upper bound is null

25 (at block 314), then the clean transform 12 inserts the lower replacement value in the field (at block 316) if the field value is less than the lower bound minus any numeric tolerance indicated in sub-parameter 100. If only the lower bound is NULL (at block 318), then the clean transform 12 inserts the upper replacement value in the field (at block 320) if

the field value is greater than the upper bound plus any numeric tolerance. If neither the lower nor upper bound are NULL, then the clean transform 12 inserts (at block 322) the lower replacement value if the field is less than or equal to the lower bound minus any numeric tolerance OR inserts the upper replacement value if the field is greater than or equal to the upper bound plus any tolerance. After applying one rule against one field at blocks 310, 316, 320 or 322, the clean transform 12 then proceeds (at block 323) back to block 307 to consider the next rule in the rule table, if there are further rules. After applying all rules in a rule table to one field (row), the clean transform 12 proceeds (at block 312) back to block 36 to consider the next field (row) in the input data column.

10

After processing all rules in the rule table and rows in the input data column from block 312 or block 324, the clean transform 12 proceeds (at block 324) to block 204 to consider any further rule definitions in the replacement rule definition parameter 60. After the last rule definition 80 in the replacement rule definitions parameter 60 is processed at blocks 220, 274 or 324, then the clean transform 12 then determines (at block 326) whether the output data table name 54 specifies an output data table 18 in the database 10 to receive any cleaned or modified input data columns. If so, then the clean transform 12 writes (at block 328) the input data columns, including replaced and cleaned fields, to the specified output data table 19. Otherwise, the modified and processed input data columns are written (at block 330) to the input table 14.

Preferred embodiments provide a command data structure to control a stored procedure program to clean columns of data in an input table the database 10. The rules to clean the tables may be provided in a rule table data structure stored in the server 4 that provides one or more clean rules for different columns and in different sequences. This allows a single instance of execution of the clean program to perform numerous types of operations on the columns in the input table. Preferred embodiments thus provide a flexible and programmable data structure and program to provide fine grained control of clean operations. Further, with preferred embodiments, the client 6 does not increase

network traffic because the client does not transfer the tables or the rule tables to clean between the database server 4. Instead, the client 6 merely provides a command data structure including various parameters and rules to a stored procedure 12 that executes in the server 4 to perform the clean operations within the database program 8 on the server

5 4. Such savings in network traffic can be significant when very large database tables, including millions or billions, of records are cleaned.

FIGs. 7 and 8a, b, c, d illustrate an example of how rule tables may be applied to clean an input data table. FIG. 7 illustrates an input data table to be cleaned including columns concerning sales information: city, indicating the sale location; product name;

10 product category; dollar sales; and unit sales. FIGs. 8a, b, c, and d illustrate different rule tables to apply to columns in the input table in FIG. 7. FIG. 8a illustrates a find and replace rule to locate certain strings in the City column that begin and end with "%". The first percent in the find value represents any string of zero or more characters. An escape character of "/" is indicated in parameter 96 so the search criteria includes strings that

15 start and end with the percent sign (%). The replace value is on the right. Thus, FIG. 8a illustrates a rule table including multiple rules having search criteria to apply against each field in a column, e.g., the CITY column, of the input data table. FIG. 8a would remove the abbreviations in the CITY input data column and replace them with a complete city name for the replacement value corresponding to the find value matching the field in the

20 City column.

FIG. 8b provides a numeric clip rule table. In the Unit Sales column, values less than 100 are considered invalid and will be rounded up to 100, and any sales number greater than 15,000 is considered NULL. Later, the rule table in FIG. 8d will be used to flag rows with a NULL value to avoid in calculations. Thus, the rule table in FIG. 8b

25 processes each field in the Unit Sales data column and replaces values less than 100 with 100 and replaces values greater than 15,000 with NULL. A numeric tolerance may be specified to apply the rule to values "close enough" to the lower or upper bound.

FIG. 8c is a discretization rule table used to replace the code values in the Product Category column of the table in FIG. 7 with meaningful descriptive terms. Code values within certain string ranges are a particular type of product. For instance, a Category code in the range from "HA" to "JZ" is a hardware product. The application of the rule table in FIG. 8c determines which rule has a lower and upper bound that includes the Product Category field and applies the corresponding Replace Value for the rule having the matching find bounds. The Upper Bound in the last row of the rule table in FIG. 8c includes a NULL value. This means the replacement value is inserted in any field in the Product Category column having a string value greater than "RX." FIG. 8c further provides a sort column, which would be included in the sort-key column name parameter 88, indicating the order in which to sort the rows in the rule table before applying the rule table to the input column

FIG. 8d is a find and replace table used to clean up the already processed data in the input table in FIG. 7. A rule definition including the rule table in FIG. 8d and having the row clean indicator parameter 94 set to YES would eliminate from further processing and from the final revised table any row in the input table (FIG. 7) having a product name of NULL. If the compress white space indicator is set in the parameters with the rule definition including FIG. 8d applied to the Product Name column, then any white space will be removed from the product name, such as the spaces between "Pick Axe" in the first data row of the input table in FIG. 7.

Conclusion

This concludes the description of the preferred embodiments of the invention. The following describes some alternative embodiments for accomplishing the present invention.

The preferred embodiments may be implemented as a method, apparatus or article of manufacture using standard programming and/or engineering techniques to produce software, firmware, hardware, or any combination thereof. The term "article of

manufacture” (or alternatively, “computer program product”) as used herein is intended to encompass one or more computer programs and data files accessible from one or more computer-readable devices, carriers, or media, such as a magnetic storage media, “floppy disk,” CD-ROM, a file server providing access to the programs via a network

5 transmission line, holographic unit, etc. Of course, those skilled in the art will recognize that many modifications may be made to this configuration without departing from the scope of the present invention.

Preferred embodiments were described with respect to specific data structures, such as a rule table having columns of rules, and an arrangement of parameters to provide
10 a vehicle for transferring commands to the clean transform stored procedure program. However, those skilled in the art will recognize that modifications may be made to the architecture of the data structures used to convey multiple clean rules and still remain within the scope of the present invention.

Preferred embodiments were described with respect to three rule types, find and
15 replace, discretization, and numeric clip. In further embodiments, other types of rules may be provided and included in the command data structure of the preferred embodiments to perform different types of clean operations known in the art.

In preferred embodiments, the clean transform program was executed in a stored procedure type program, such as that used in the IBM DB2 database system. However, in
20 further embodiments, different types of application programs, other than stored procedure programs, may be executed in the server 4 or even the client 6 to perform clean operations in accordance with the command data structures of the preferred embodiments.

In preferred embodiments, the input table and output table were included in a database in the server in which the clean transform program is executing. In alternative
25 embodiments, the rule, input, and output tables may be distributed at different storage locations at different network devices.

In preferred embodiments, a client constructed the clean operation command and communicated such commands to the database server. In alternative embodiments, the

65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100

clean operation command of the preferred embodiments may be executed on the machine used to construct the command.

In summary, preferred embodiments disclose a method, system, program, and data structure for performing a clean operation on an input table. The input table to clean is indicated in an input data table name. At least two one rule definition is processed to clean the input table. Each rule definition indicates a find criteria, a replacement value, and an input data column in the input table. For each rule definition, the input data column is searched for any fields that match the find criteria. The replacement value for the particular rule definition is inserted in the fields in the input data column that match the find criteria. Subsequent applications of additional rule definitions applied to the same input data column operate on replacement values inserted in the input data column during previously applied rule definitions.

The foregoing description of the preferred embodiments of the invention has been presented for the purposes of illustration and description. It is not intended to be exhaustive or to limit the invention to the precise form disclosed. Many modifications and variations are possible in light of the above teaching. It is intended that the scope of the invention be limited not by this detailed description, but rather by the claims appended hereto. The above specification, examples and data provide a complete description of the manufacture and use of the composition of the invention. Since many embodiments of the invention can be made without departing from the spirit and scope of the invention, the invention resides in the claims hereinafter appended.

**Microsoft, Windows, Windows NT are registered trademarks and Microsoft SQL Server is a trademark of Microsoft Corporation; DB2, AIX, OS/390, OS/400, and OS/2 are registered trademarks of IBM; and Oracle8 is a trademark of Oracle Corporation; and Solaris is a trademark of Sun Microsystems, Inc.

WHAT IS CLAIMED IS:

1 1. A method for performing a clean operation on an input table having an
2 input table name, comprising:
3 receiving at least one rule definition, wherein each rule definition indicates a find
4 criteria, a replacement value, and an input data column in the input table;
5 searching, for each rule definition, the input data column for any fields that match
6 the find criteria; and
7 inserting, for each rule definition, the replacement value in the fields in the input
8 data column that match the find criteria, wherein subsequent applications of additional
9 rule definitions applied to the same input data column operate on replacement values
10 inserted in the input data column in previously applied rule definitions.

1 2. The method of claim 1, wherein each rule definition is associated with one
2 rule table including the find criteria and replacement value, wherein a rule table column
3 parameter for each rule definition indicates the columns in the rule table including the
4 find criteria and replacement value for the rule definition.

1 3. The method of claim 1, wherein there is a separate rule table including the
2 find criteria and replacement value associated with at least one rule definition, wherein,
3 for each rule definition, a rule table column parameter indicates the columns in the rule
4 table for the rule definition including the find criteria and replacement value for that rule
5 definition.

1 4. The method of claim 1, wherein the input data column for a first and
2 second applied rule definitions is the same input data column, wherein the replacement
3 value for the first rule definition is inserted into at least one field in the input data
4 column, and wherein the find criteria of the second rule definition is applied to the
5 replacement value inserted in the input data column.

694220044636660

1 9. The method of claim 8, wherein the at least one rule definition including
2 find criteria having upper and lower bounds includes multiple find criteria and a
3 corresponding replacement value for each find criteria, wherein the step of searching the
4 input data column comprises applying each of the multiple find criteria to one field until
one of: (i) a match occurs and (ii) none of the multiple find criteria are found to match the
field content, and wherein inserting the replacement value comprises inserting the
replacement value corresponding to one find criteria that matched the field content.

1 14. A system for performing a clean operation on an input table having an
2 input data table name, comprising;

3 means for receiving at least one rule definition, wherein each rule definition
4 indicates a find criteria, a replacement value, and an input data column in the input table;
5 means for searching, for each rule definition, the input data column for any fields
6 that match the find criteria; and
7 means for inserting, for each rule definition, the replacement value in the fields in
8 the input data column that match the find criteria, wherein subsequent applications of
9 additional rule definitions applied to the same input data column operate on replacement
10 values inserted in the input data column in previously applied rule definitions.

1 15. The system of claim 14, wherein each rule definition is associated with
2 one rule table including the find criteria and replacement value, wherein a rule table
3 column parameter for each rule definition indicates the columns in the rule table
4 including the find criteria and replacement value for the rule definition.

1 16. The system of claim 14, wherein there is a separate rule table including the
2 find criteria and replacement value associated with at least one rule definition, wherein,
3 for each rule definition, a rule table column parameter indicates the columns in the rule
4 table for the rule definition including the find criteria and replacement value for that rule
5 definition.

1 17. The system of claim 14, wherein the input data column for a first and
2 second applied rule definitions is the same input data column, further comprising:
3 means for inserting the replacement value for the first rule definition into at least
4 one field in the input data column; and
5 means for applying the find criteria of the second rule definition to the
6 replacement value inserted in the input data column.

03126014636650

1 22. The system of claim 21, wherein the at least one rule definition including
2 find criteria having upper and lower bounds includes multiple find criteria and a
3 corresponding replacement value for each find criteria, wherein the means for searching
4 the input data column comprises applying each of the multiple find criteria to one field
5 until one of: (i) a match occurs and (ii) none of the multiple find criteria are found to
6 match the field content, and wherein the means for inserting the replacement value

7 comprises inserting the replacement value corresponding to one find criteria that matched
8 the field content.

1 23. The system of claim 20, wherein the means for searching comprises
2 searching for any fields that have values outside of one of the upper and lower bounds.

1 24. The system of claim 14, wherein the find criteria for at least one rule
2 definition comprises an upper bound and lower bound and wherein the replacement value
3 is an upper replacement value and further comprising a lower replacement value, wherein
4 the means for searching comprises searching for any fields that have values within the
5 upper and lower bounds and wherein inserting comprises inserting the upper replacement
6 value if the field has a value greater than the upper bound and inserting the lower
7 replacement value if the field has a value less than the lower bound.

1 25. The system of claim 24, wherein the at least one rule definition including
2 find criteria having upper and lower bounds includes multiple find criteria and a
3 corresponding upper and lower replacement value for each find criteria, wherein the
4 means for searching the input data column comprises applying each of the multiple find
5 criteria to one field until one of: (i) a match occurs and (ii) none of the multiple find
6 criteria are found to match the field content, and wherein the means for inserting the
7 replacement value comprises inserting the replacement value corresponding to one find
8 criteria that matched the field content.

1 26. The system of claim 14, wherein the rule definitions include a row clean
2 flag, and wherein at least one rule definition has the row clean flag set, further comprising
3 removing any row including a field matching the search criteria from the input table when
4 the row clean flag is set.

691260-463660

9 inserting, for each rule definition, the replacement value in the fields in the input
10 data column that match the find criteria, wherein subsequent applications of additional
11 rule definitions applied to the same input data column operate on replacement values
12 inserted in the input data column in previously applied rule definitions.

1 28. The article of manufacture of claim 27, wherein each rule definition is
2 associated with one rule table including the find criteria and replacement value, wherein a
3 rule table column parameter for each rule definition indicates the columns in the rule
4 table including the find criteria and replacement value for the rule definition.

1 29. The article of manufacture of claim 27, wherein there is a separate rule
2 table including the find criteria and replacement value associated with at least one rule
3 definition, wherein, for each rule definition, a rule table column parameter indicates the
4 columns in the rule table for the rule definition including the find criteria and replacement
5 value for that rule definition.

1 30. The article of manufacture of claim 27, wherein the input data column for
2 a first and second applied rule definitions is the same input data column, wherein the
3 replacement value for the first rule definition is inserted into at least one field in the input

[illegible]

1 data column, and wherein the find criteria of the second rule definition is applied to the
2 replacement value inserted in the input data column.

1 31. The article of manufacture of claim 27, wherein at least one rule definition
2 includes multiple find criteria and a corresponding replacement value for each find
3 criteria, wherein the step of searching the input data column comprises applying each of
4 the multiple find criteria to one field until one of: (i) a match occurs and (ii) none of the
5 multiple find criteria are found to match the field content, and wherein inserting the
6 replacement value comprises inserting the replacement value corresponding to one find
7 criteria that matched the field content.

1 32. The article of manufacture of claim 31, wherein a sort column includes
2 values to use to sort the multiple find criteria and corresponding replacement value,
3 wherein the step of searching comprises applying the multiple find criteria to each field in
4 the order specified in the sort column.

1 33. The article of manufacture of claim 27, wherein the rule definition
2 comprises a type of rule that is a member of the set of rules consisting of: find and
3 replace, discretization, and numeric clip, wherein at least two rule definitions are
4 comprised of different rule types.

1 34. The article of manufacture of claim 27, wherein the find criteria for at least
2 one rule definition comprises an upper bound and lower bound, wherein searching
3 comprises searching for any fields that have values within the upper and lower bounds.

1 35. The article of manufacture of claim 34, wherein the at least one rule
2 definition including find criteria having upper and lower bounds includes multiple find
3 criteria and a corresponding replacement value for each find criteria, wherein the step of

09120043560

1 36. The article of manufacture of claim 34, wherein searching comprises
2 searching for any fields that have values outside of one of the upper and lower bounds.

37. The article of manufacture of claim 27, wherein the find criteria for at least one rule definition comprises an upper bound and lower bound and wherein the replacement value is an upper replacement value and further comprising a lower replacement value, wherein searching comprises searching for any fields that have values within the upper and lower bounds and wherein inserting comprises inserting the upper replacement value if the field has a value greater than the upper bound and inserting the lower replacement value if the field has a value less than the lower bound.

1 38. The article of manufacture of claim 37, wherein the at least one rule
2 definition including find criteria having upper and lower bounds includes multiple find
3 criteria and a corresponding upper and lower replacement value for each find criteria,
4 wherein the step of searching the input data column comprises applying each of the
5 multiple find criteria to one field until one of: (i) a match occurs and (ii) none of the
6 multiple find criteria are found to match the field content, and wherein inserting the
7 replacement value comprises inserting the replacement value corresponding to one find
8 criteria that matched the field content.

1 39. The article of manufacture of claim 27, wherein the rule definitions
2 include a row clean flag, and wherein at least one rule definition has the row clean flag

[illegible]

[illegible]

1 41. The memory device of claim 40, wherein at least one rule definition
2 further includes:
3 indication of one rule table including the find criteria and replacement value for
4 the at least two rule definitions, such that the one rule table includes the find criteria and
5 replacement value for the at least two rule definitions; and
6 a rule table column parameter for the at least two rule definitions indicating the
7 columns in the rule table including the find criteria and replacement value for the rule
8 definitions.

1 42. The memory device of claim 40, wherein at least one rule definition
2 further includes:
3 indication of a separate rule table for each rule definition including the find
4 criteria and replacement value for the at least two rule definitions; and

1 a rule table column parameter indicating the columns in the rule table for the rule
2 definition including the find criteria and replacement value for that rule definition.

1 43. The memory device of claim 40, wherein the input data column for a first
2 and second applied rule definitions is the same input data column.

1 44. The memory device of claim 40, wherein at least one rule definition
2 further includes:
3 multiple find criteria and a corresponding replacement value for each find criteria,
4 wherein the input data column is searched by applying each of the multiple find criteria to
5 one field until one of: (i) a match occurs and (ii) none of the multiple find criteria are
6 found to match the field content, and wherein the replacement value corresponding to the
7 matching find criteria is inserted into the field.

1 45. The memory device of claim 4, wherein the at least one rule definition
2 further comprises a sort column including values to use to sort the multiple find criteria
3 and corresponding replacement value, wherein the multiple find criteria are applied to
4 each field in the input data column in the order specified in the sort column.

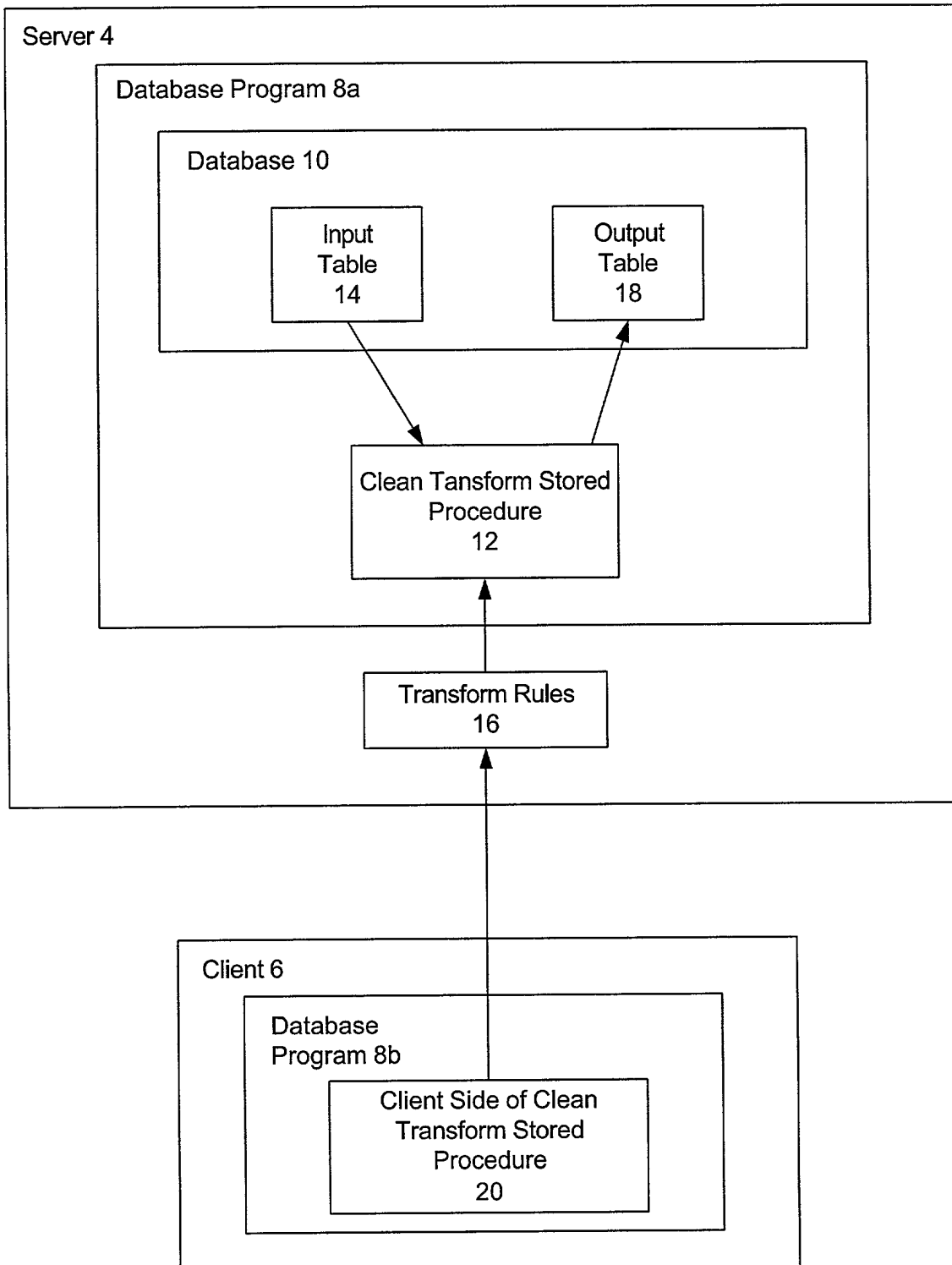
1 46. The memory device of claim 40, wherein the rule definition comprises a
2 type of rule that is a member of the set of rules consisting of: find and replace,
3 discretization, and numeric clip, wherein at least two rule definitions are comprised of
4 different rule types.

ABSTRACT

[illegible]

FIG. 1

2



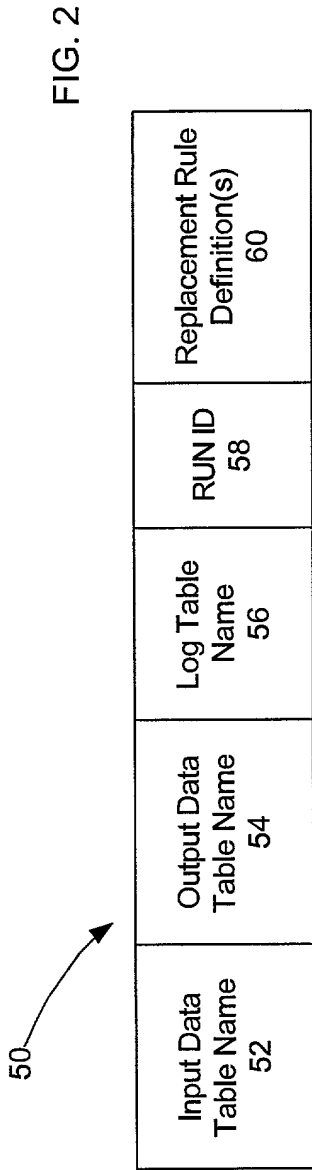


FIG. 3

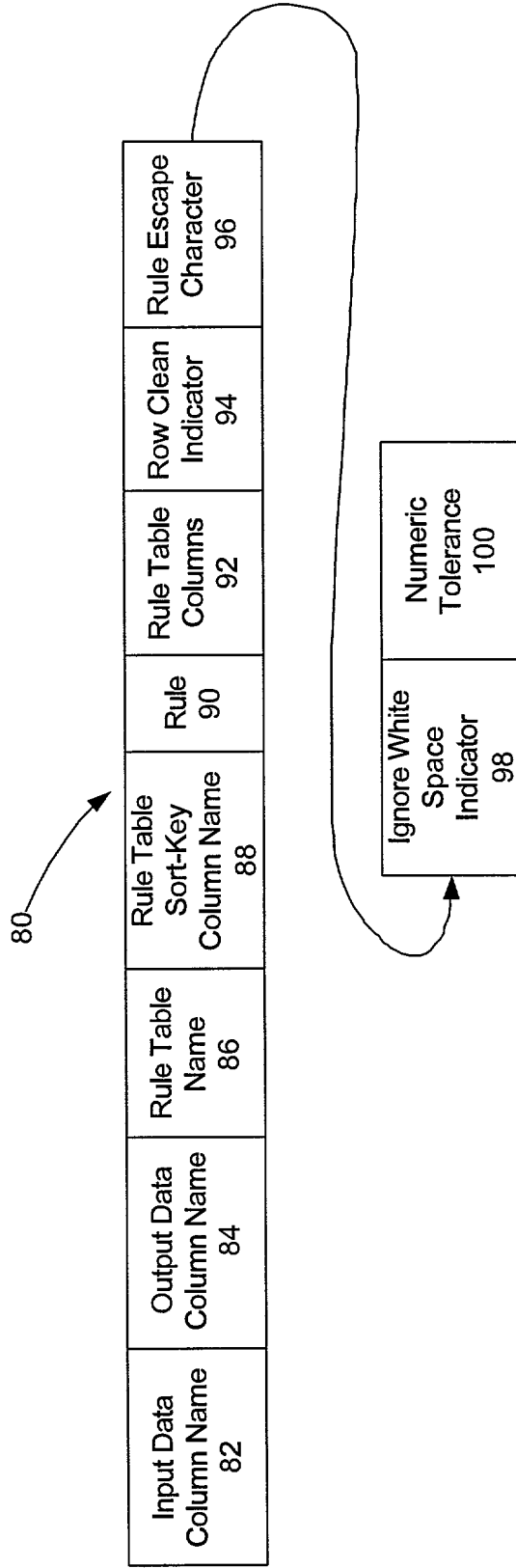


FIG. 3a

find pattern	replacement text value
--------------	------------------------

FIG. 3b

find value	replacement text value
------------	------------------------

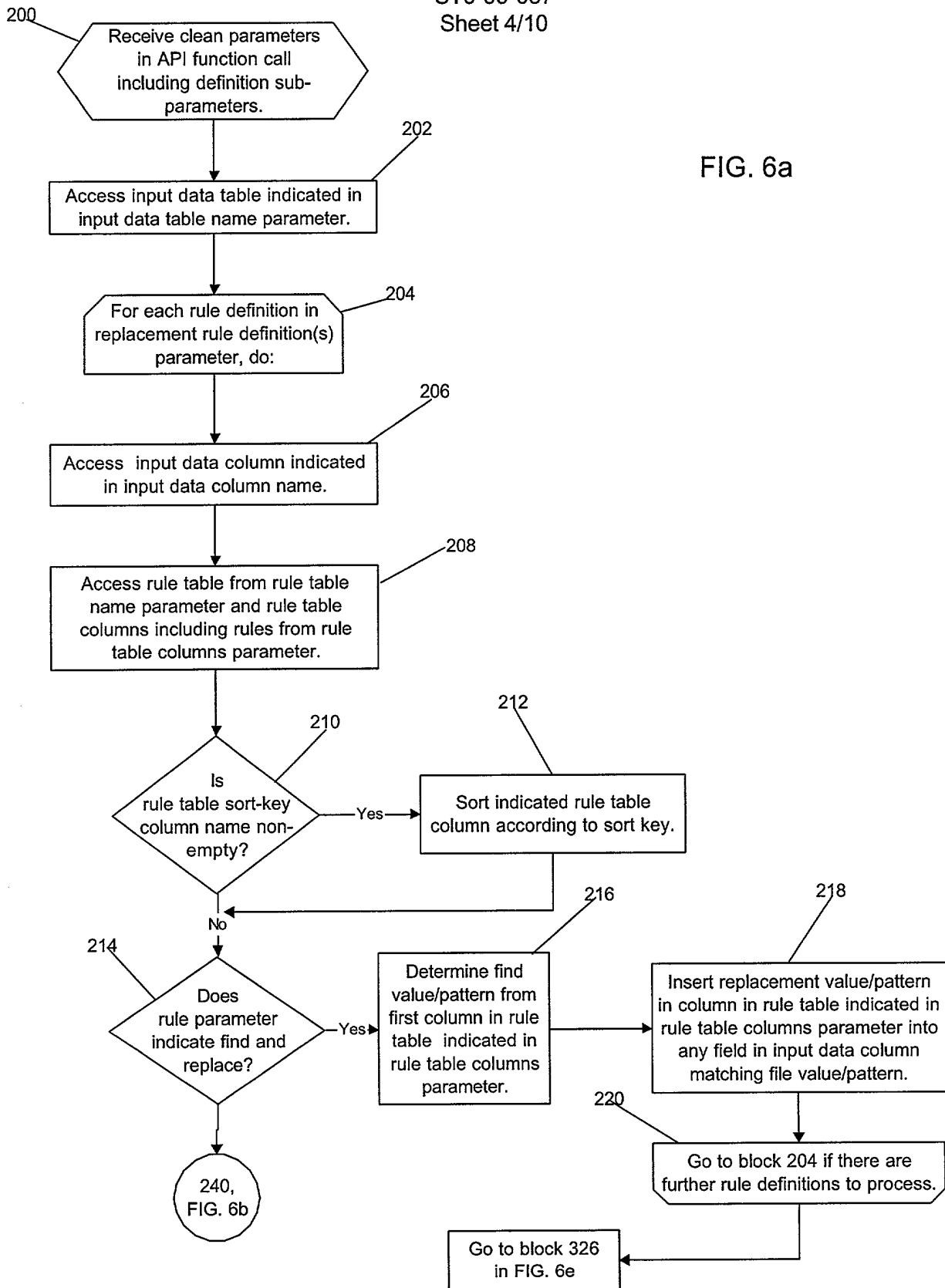
FIG. 4

lower bound	upper bound	replacement value
-------------	-------------	-------------------

FIG. 5

lower bound	lower replacement value	upper bound	upper replace value
-------------	-------------------------	-------------	---------------------

FIG. 6a



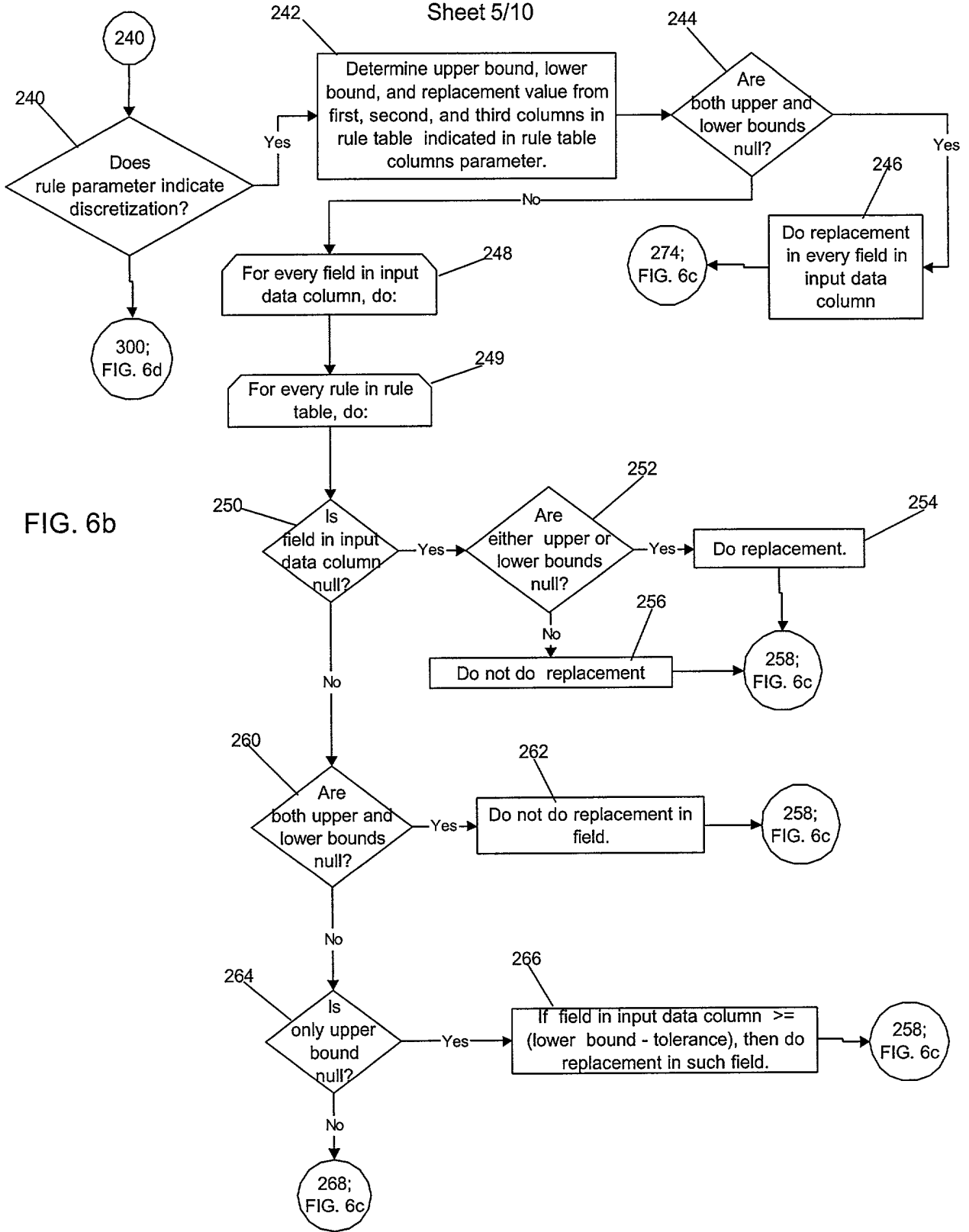
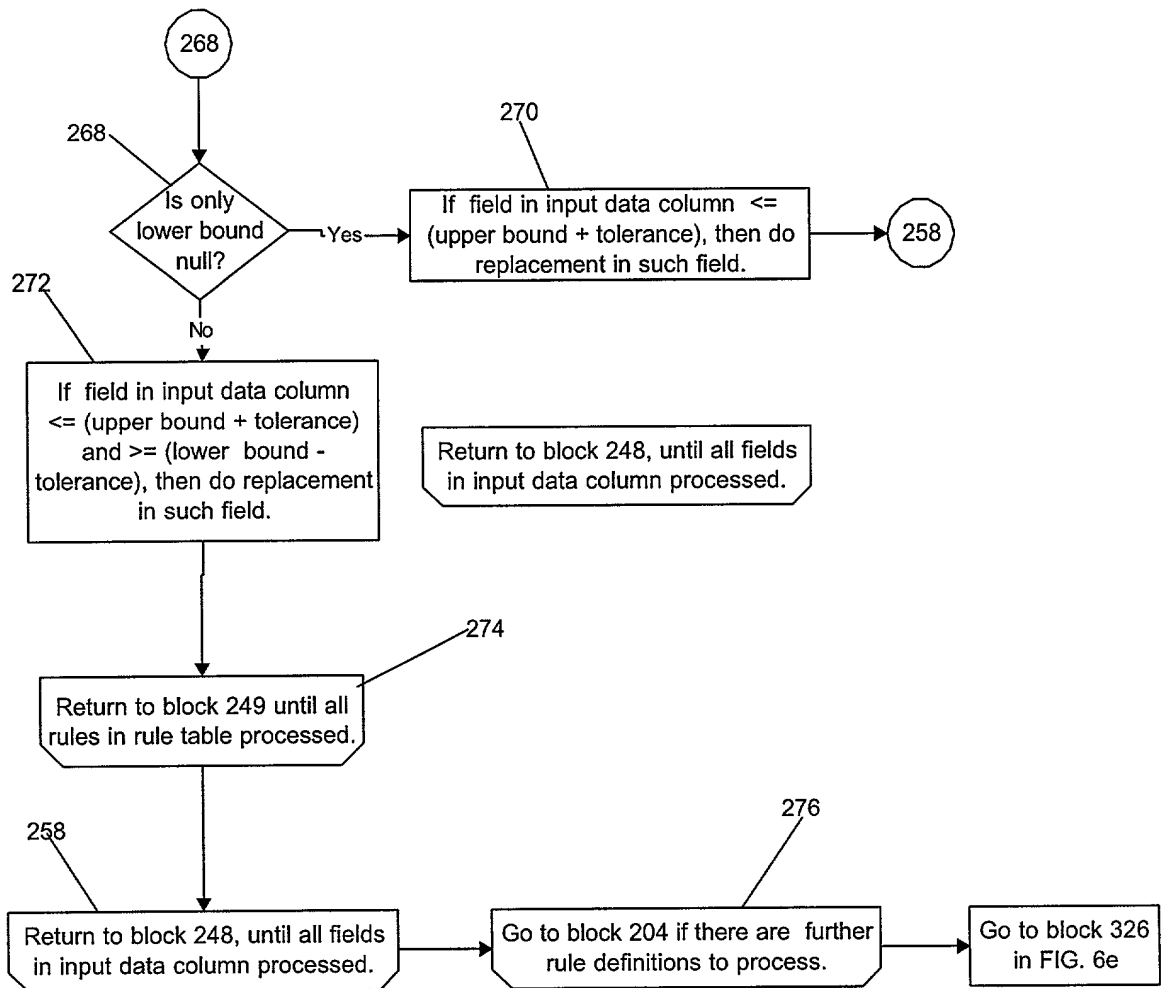


FIG. 6c



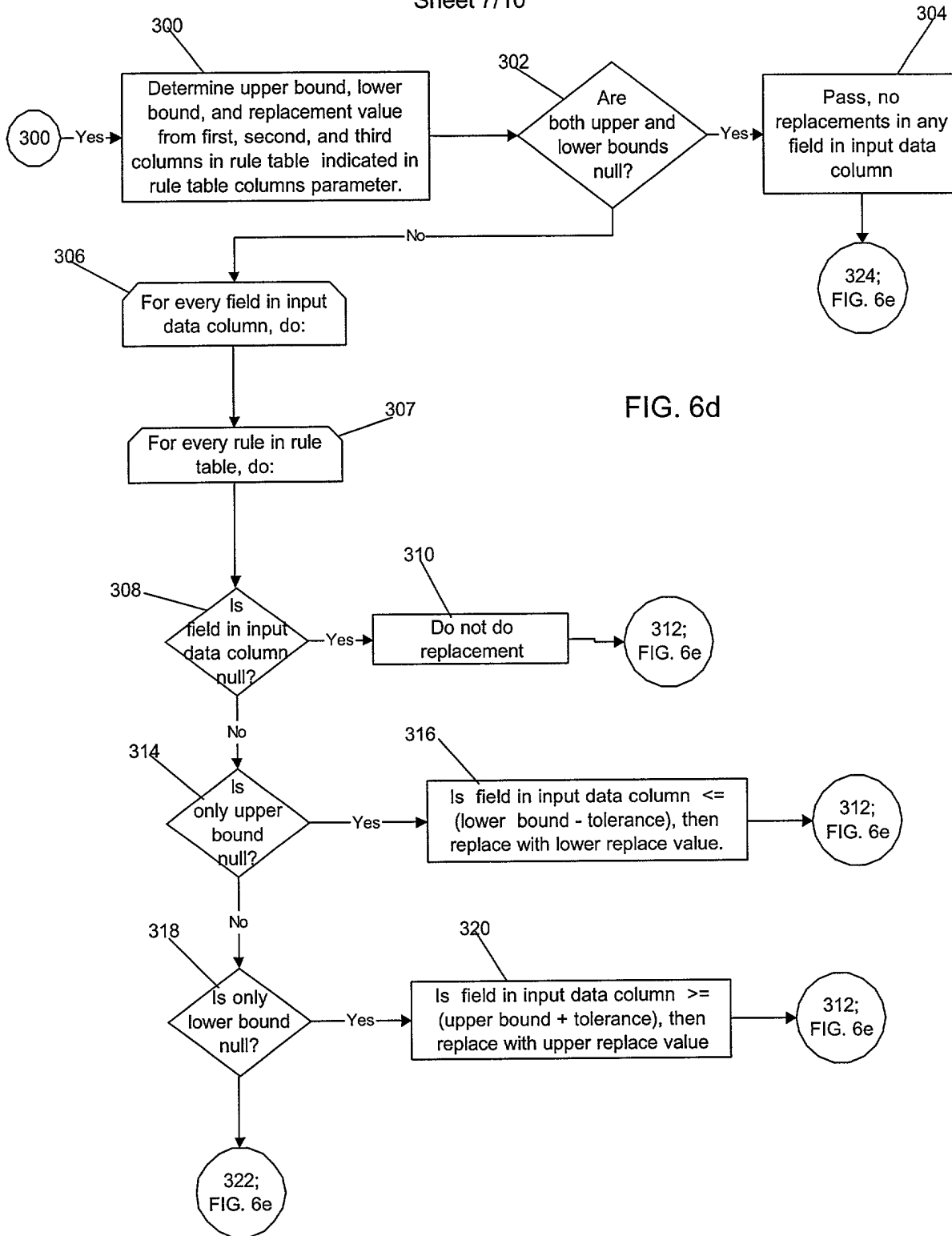


FIG. 6e

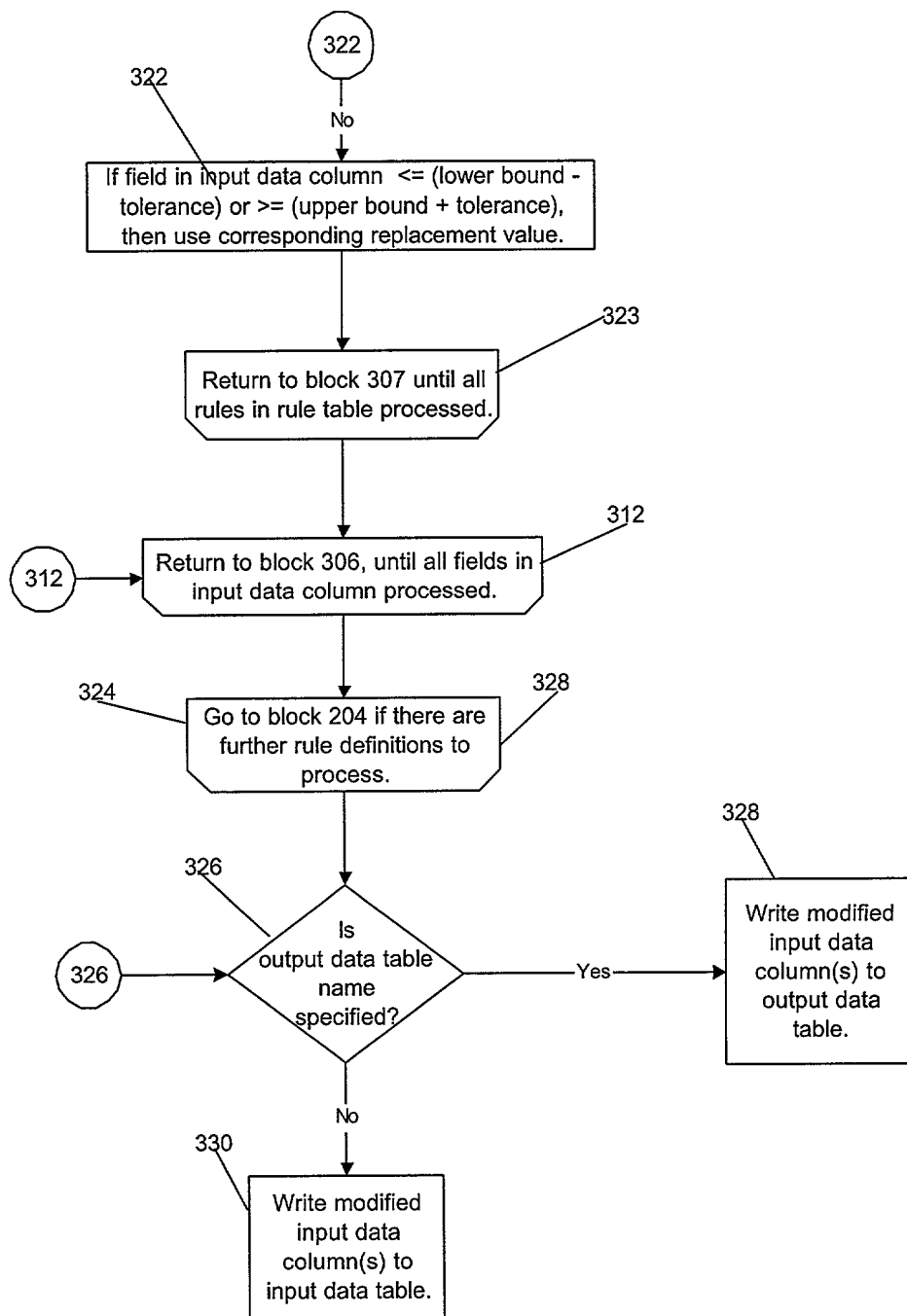


FIG. 7

CITY	Product Name	Product Category	Sales (\$)	Unit Sales
Varchar (20)	Varchar(30)	Varchar (15)	Decimal(8,2)	Integer
11%LAX455	Pick Axe	IATE	5897.32	1100
SJC1298%SEA	Small Shovel	JXJ8972	456.98	43
%lax6%SFO94	Tarp	HAU1245	893.12	422
%SJC34LAX45	Flares	AX111	1200.45	1189
4515%SEA	NULL	AC183	333.9	23
3452%SJC34	""	MIS222	1818.4	15002

FIG. 8a

Find Value	Replace Value
Varchar(10)	Varchar(20)
%%LAX%	Los Angeles
%%DFW%	Dallas/Fort Worth
%%SJC%	San Jose, CA
%%SFO%	San Francisco
%%SEA%	Seattle/Tacoma

FIG. 8b

Lower Bound	Lower Replacement	Upper Bound	Upper Replacement
Integer	Integer	Integer	Integer
100	100	15000	NULL

FIG. 8c

Sort Column	Lower Bound	Upper Bound	Replace Value
Integer	Varchar(6)	Varchar(6)	Varchar(10)
1	I	IZZZZ	Special
2	HA	JZ	Hardware
3	LA	NZ	Miscel
4	AA	BJ	Winter Storm
5	RX	NULL	Winter Storm

FIG. 8d

Find Value	Replace Value
Varchar(5)	Varchar(10)
NULL	""

2025 RELEASE UNDER E.O. 14176

As a below named inventor, I hereby declare that:

My residence, post office address and citizenship are as stated below next to my name;

I believe I am the original, first and sole inventor (if only one name is listed below) or an original, first and joint inventor (if plural names are listed below) of the subject matter which is claimed and for which a patent is sought on the invention entitled

METHOD, SYSTEM, PROGRAM, AND DATA STRUCTURE FOR CLEANING A DATABASE TABLE

the specification of which (check one)

☒ is attached hereto.

☐ was filed on _____

as Application Serial No. _____

and was amended on _____ (if applicable).

I hereby state that I have reviewed and understand the contents of the above identified specification, including the claims, as amended by any amendment referred to above.

I acknowledge the duty to disclose information which is material to patentability as defined in Title 37, Code of Federal Regulations, Section 1.56.

I hereby claim foreign priority benefits under Title 35, United States Code, Section 119(a)-(d) or Section 365(b) of any foreign application(s) for patent or inventor's certificate, or Section 365(a) of any PCT International application which designated at least one country other than the United States, listed below and have also identified below any foreign application for patent or inventor's certificate or PCT International application having a filing date before that of the application on which priority is claimed:

Prior Foreign Application(s)

Priority Claimed

☐ None ☐ Yes ☐ No
(Number) (Country) (Day/Month/Year Filed)

I hereby claim the benefit under Title 35, United States Code, Section 120 of any United States application(s) or Section 365(c) of any PCT International application designating the United States, listed below and, insofar as the subject matter of each of the claims of this application is not disclosed in the prior United States or PCT International application in the manner provided by the first paragraph of Title 35, United States Code, Section 112, I acknowledge the duty to disclose information which is material to patentability as defined in Title 37, Code of Federal Regulations, Section 1.56, which occurred between the filing date of the prior application and the national or PCT international filing date of this application:

☐ None
(Application Serial No.) (Filing Date) (Status) (patented, pending, abandoned)

I hereby declare that all statements made herein of my own knowledge are true and that all statements made on information and belief are believed to be true; and further that these statements were made with the knowledge that willful false statements and the like so made are punishable by fine or imprisonment, or both, under Section 1001 of Title 18 of the United States Code and that such willful false statements may jeopardize the validity of the application or any patent issued thereon.

POWER OF ATTORNEY: As a named inventor, I hereby appoint the following attorney(s) and/or agent(s) to prosecute this application and transact all business in the Patent and Trademark Office connected therewith. (list name and registration number) Romualdas Strimaitis, Reg. No. 35,697; Prentiss W. Johnson, Reg. No. 33,123; Ingrid M. Foerster, Reg. No. 36,511; Timothy M. Farrell, Reg., No. 37,321; Christopher A. Hughes, Reg. No. 26,914; John E. Hoel, Reg. No. 26,279; Edward A. Pennington, Reg. No. 32,588; Joseph C. Redmond, Jr., Reg. No. 18,753; David W. Victor, Reg. No. 39,867; William K. Konrad, Reg. No. 28,868; Alan S. Raynes, Reg. No. 39,809.

Send correspondence to:

David Victor, Esq
Konrad Raynes & Victor, LLP
1180 South Beverly Dr., Ste. 501
Los Angeles, CA 90035

Direct all telephone calls to David Victor at (310) 553-7977

FULL NAME OF INVENTOR ONE: Mark Anthony Cesare

INVENTORS SIGNATURE: 

DATE:

Sept 16, 1999

RESIDENCE: 243 Mangels Avenue, San Francisco, California 94131

CITIZENSHIP: United States

POST OFFICE ADDRESS: same as residence

FULL NAME OF INVENTOR TWO: Tom Robert Christopher

INVENTORS SIGNATURE: 

DATE:

Sept. 16, 1999

RESIDENCE: 18315 Murphy Springs Dr., Morgan Hill, California 95037

CITIZENSHIP: United States

POST OFFICE ADDRESS: same as residence

FULL NAME OF INVENTOR THREE: Julie Ann Jerves

INVENTORS SIGNATURE: 

DATE:

Sept 15, 1999

RESIDENCE: 20099 Sea Gull Way, Saratoga, California 95070

CITIZENSHIP: United States

POST OFFICE ADDRESS: same as residence

